

Reconstructing the distribution of haloes and mock galaxies below the resolution limit in cosmological simulations

Sylvain de la Torre^{*} and John A. Peacock

SUPA[†], Institute for Astronomy, University of Edinburgh, Royal Observatory, Blackford Hill, Edinburgh EH9 3HJ, UK

ABSTRACT

We present a method for populating dark matter simulations with haloes of mass below the resolution limit. It is based on stochastically sampling a field derived from the density field of the halo catalogue, using constraints from the conditional halo mass function $n(m|\delta)$. We test the accuracy of the method and show its application in the context of building mock galaxy samples. We find that this technique allows precise reproduction of the two-point statistics of galaxies in mock samples constructed with this method. Our results demonstrate that the full information content of a simulation can be communicated efficiently using only a catalogue of the more massive haloes.

Key words: Cosmology: large-scale structure of Universe – Galaxies: statistics.

1 INTRODUCTION

The distribution of luminous matter in the Universe is a rich source of information concerning the large-scale mass distribution, the formation and evolution of baryonic structures, plus the overall properties of the Universe. In this way, future large galaxy surveys will yield extremely precise measurements of quantities such as the global expansion history and growth rate of structure from measurements of the galaxy power spectrum (e.g. [Laureijs et al. 2011](#); [Schlegel et al. 2011](#)). But achieving high statistical precision is only possible if we have a full understanding of all the potential systematic biases arising from survey selection – and realistic surveys are sufficiently complex that the necessary calibration of statistical methods can only be achieved by using mock survey realisations. Moreover, mock datasets represent one of the most efficient ways of estimating uncertainties related to the statistics of interest, including both estimation and sample variance errors; typically 100-1000 independent realisations are required for this purpose. Naturally, robust answers to these questions require the mocks to be as realistic as possible, although they do not need to match reality in every respect, so long as they contain the main complicating factors that could be a source of error in the real data.

For these reasons, dark matter N-body simulations coupled with semi-analytical treatments of galaxy formation or halo occupation distribution (HOD) techniques have become standard methods for producing mock survey samples. These two-layered methods reflect our current understanding of cosmology and galaxy formation, where the large-scale structure of the Universe, dominated by dark matter, evolves through gravity; galaxies then form inside the dark matter haloes by the collapse and cooling of baryonic gas into the potential wells that they provide. However, dark matter N-

body simulations are finite and their usefulness is limited by the volume and mass resolution that they can probe. For cosmological surveys, the probed volume directly influences the statistical error with which one can measure the cosmological parameters. Conversely, mass resolution is more important for galaxy evolution surveys, in which one is more interested in a complete census of galaxies, i.e. probing the whole range of galaxy masses and associated physical properties. These competing criteria of volume and resolution mean that all existing simulations are a compromise.

One way of tackling this limitation is to make an approximate reconstruction of the distribution of low-mass haloes below the resolution limit. Given predictions of the halo bias and abundance at the lowest masses, one can aim to recover the missing information below the resolution limit. In principle, a good deal of information is available about the distribution of low-mass haloes, since about half of the mass in most simulations is not resolved into haloes. But for large simulation volumes, handling this quantity of particle data is cumbersome, and often only the catalogue of resolved haloes can be efficiently transmitted. Our aim here is thus to see to what extent the distribution of all haloes can be inferred from only those that are detected in a simulation. We present in this letter a method of this sort, which we show to be particularly effective in the context of building mock galaxy surveys.

2 METHOD

The proposed method consists of two steps: one first uses the simulation halo catalogue to estimate a halo density field which is then sampled stochastically to obtain haloes with mass below the resolution limit and with the correct abundance and bias. To predict the number of haloes of different masses in each region of the simulation, i.e. the conditional halo mass function, we use the peak-background split formalism ([Bardeen et al. 1986](#); [Cole & Kaiser](#)

^{*} E-mail: sdlt@roe.ac.uk

[†] Scottish Universities Physics Alliance

1989). The shape the halo mass function $n(m)$ and bias factor $b(m)$, which are the basic ingredients entering in the conditional halo mass function, have to be extrapolated to the masses below the nominal minimum halo mass of the simulation or to be assumed from theory. In the following subsections we describe in detail the two parts of the method.

2.1 Halo density field estimation

The main idea is to use the simulation halo catalogue, which preferentially traces the densest environments, to infer the full range of overdensity. The first step is to estimate the halo density field traced by the haloes originally present in the simulation. There are several ways to estimate the density field, the simplest being to count the number of objects on a cubical grid, associating each halo to the closest grid node. Generally the means of assigning objects to grid nodes and the grid size have an impact on the accuracy of the recovered density field. Optimally, we would like to use cells as small as possible, so as to probe the smallest-scale density fluctuations, but also large enough to avoid introducing shot noise. The optimal grid size will then depend on the number density of haloes in the simulation and, in turn, on the nominal halo mass resolution. One way of reducing the shot noise in the reconstructed density field is to use Delaunay tessellation. In that case, instead of using fixed-size cells to estimate number densities, one uses tetrahedra whose size varies adaptively depending on the local number of objects. The resulting density field estimates can then be interpolated onto a fine grid for convenience. This method allows the reduction of the shot noise contribution, while retaining high-resolution information when it is available. Other adaptive smoothing methods based on e.g. nearest neighbours could also be used. We show in the next section the improvement on the halo density field estimation that this can produce.

2.2 Low-mass halo population

Once a continuous halo density field is estimated, one can use the expected number of haloes of mass m in each cell of mass overdensity δ , i.e. the conditional halo mass function $n(m|\delta)$, to populate the simulation with haloes of mass below the resolution limit. The halo density field δ_h is biased with respect to the mass density field δ and consequently has to be de-biased prior to being used to predict the number of expected low-mass haloes.

We follow the peak-background split formalism and write the conditional halo mass function as,

$$n(m|\delta) = n(m)(1 + \langle \delta_h(m) | \delta \rangle), \quad (1)$$

where $n(m)$ is the (unconditional) halo mass function and $\langle \delta_h(m) | \delta \rangle$ is the function describing the biasing of haloes of mass m . In the case of sufficiently large cells, density fluctuations become linear and we can assume $\delta_h = b(m)\delta$. In this limit Equation (1) simplifies to,

$$n(m|\delta) = n(m)(1 + b(m)\delta) \quad (2)$$

where $b(m)$ is the large-scale linear halo bias factor. In practice, $n(m)$ and $b(m)$ have to be specified for mass values below m_{lim} , the minimum halo mass of the simulation. For this one can either use analytical forms or extrapolate these functions in the simulation itself. The extrapolation is relatively straightforward because those functions show only weak and relatively easily predictable variations with halo mass in the low-mass regime.

Equation (2) is valid for densities estimated on large scales where non-linear fluctuations are smeared out. However we would like to have a model that accounts to some extent for bias non-linearities which are present on small scales. One simple (local) non-linear biasing model that we can use is the power-law bias model (e.g. Mann et al. 1998; Narayanan et al. 2000) for which the halo bias is defined as,

$$1 + \delta_h \propto (1 + \delta)^{b(m)}. \quad (3)$$

This model has a certain number of advantages: it naturally avoids negative densities and depends only on one parameter. Furthermore such a power-law model has empirical support to the extent that it gives a good match to the relative biasing of different classes of galaxies (Wild et al. 2005). We will show in Section 3 that it is accurate enough for the purpose of the present method. While using the power-law bias model in Equation (1), one obtains a conditional halo mass function of the form,

$$n(m|\delta) \propto n(m)(1 + \delta)^{b(m)}. \quad (4)$$

Because the halo density field is biased and the mass overdensity that enters in Equation (4) is unknown *a priori*, one has to rewrite the conditional mass function in terms of the halo overdensity δ_h . If we assume the same biasing model to de-bias the original halo density field, then the final conditional halo mass function that we can use to populate the simulation in low-mass haloes is,

$$n(m|\delta_h) \propto n(m)(1 + \delta_h)^{b(m)/b_0}, \quad (5)$$

where b_0 is the effective bias of the original halo population defined as,

$$b_0 = \int_{m_{\text{lim}}}^{\infty} b(m)n(m) dm. \quad (6)$$

Note that there would be a factor m inside the integral if we had chosen to weight haloes by mass. But number weighting reduces both non-linear bias and shot noise from finite numbers of haloes. In practice the normalisation of Equation (5) is imposed empirically by requiring $\langle \delta_h \rangle = 0$ when volume averaging over all cells of the simulation. Finally, the number of low-mass haloes in each cell is randomly drawn by Poisson sampling the $n(m|\delta_h)$. We note that with this procedure, the low-mass haloes will not exhibit any clustering on scales below the size of the cells.

It is worth repeating that one could of course have used the dark matter particles in the simulation and worked directly from the mass density field δ using Equation (4), but using the halo catalogue provides a much more efficient means of accessing this information. In particular, it allows the reconstruction method to be applied to public simulation datasets where the full particle distribution is typically not made available.

3 TESTS ON SIMULATION DATA

We test the reconstruction method on the Millennium simulation (Springel et al. 2005) which probes a volume of $0.125 h^{-3} \text{ Gpc}^3$ with a mass resolution of $m_p = 8.6 \times 10^8 h^{-1} M_\odot$ in a Λ CDM cosmology with $(\Omega_m, \Omega_\Lambda, \Omega_b, h, n, \sigma_8) = (0.25, 0.75, 0.045, 0.7, 0.95, 0.9)$. We will also make use in the following of the MultiDark Run 1 (MDR1) dark matter N-body simulation (Prada et al. 2012). MDR1 probes a larger volume of $1 h^{-3} \text{ Gpc}^3$ with a mass resolution of $m_p = 8.721 \times 10^9 h^{-1} M_\odot$ in a Λ CDM cosmology with $(\Omega_m, \Omega_\Lambda, \Omega_b, h, n, \sigma_8) = (0.27, 0.73, 0.0469, 0.7, 0.95, 0.82)$. In both simulations, the

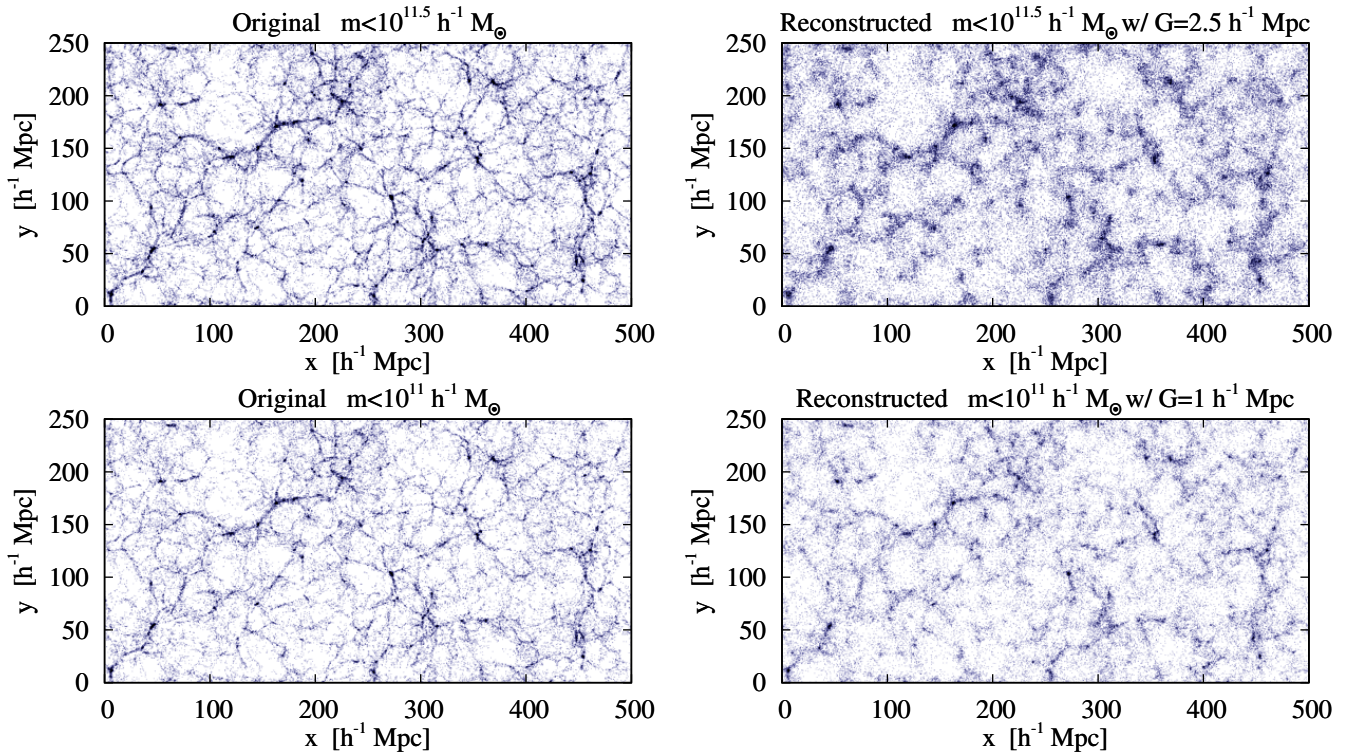


Figure 1. Comparison of the continuous density fields of original (left panels) and reconstructed haloes (right panels) in a slice of $500 \times 250 \times 15 h^{-3} \text{ Mpc}^3$ from the Millennium simulation, for two cuts in halo mass corresponding to $m < 10^{11.5} h^{-1} M_{\odot}$ (top panels) and $m < 10^{11} h^{-1} M_{\odot}$ (bottom panels). In the $m < 10^{11.5} h^{-1} M_{\odot}$ case, the reconstruction used a grid of size $G = 2.5 h^{-1} \text{ Mpc}$, while in the $m < 10^{11} h^{-1} M_{\odot}$ case, a grid of size $G = 1 h^{-1} \text{ Mpc}$ was used.

dark matter haloes have been identified from the dark matter particle distribution using a friends-of-friends algorithm and we use only the haloes identified in the snapshots at $z = 0.1$. The minimum halo mass in the Millennium and MultiDark halo catalogues are respectively $m_{\text{lim}} = 10^{10.5} h^{-1} M_{\odot}$ and $m_{\text{lim}} = 10^{11.5} h^{-1} M_{\odot}$.

We estimate the halo density field by measuring the halo density contrast defined as $\delta_h(\mathbf{r}) = (N(\mathbf{r}) - \langle N \rangle) / \langle N \rangle$ where $N(\mathbf{r})$ and $\langle N \rangle$ are respectively the number of haloes in a cell centred at position \mathbf{r} and the mean number of haloes per cell. Given the halo number density, the optimal choice of cell size falls between $2.5 h^{-1} \text{ Mpc}$ and $5 h^{-1} \text{ Mpc}$, so to have a few haloes per cell on average. We choose a grid size of $G = 2.5 h^{-1} \text{ Mpc}$ and estimate the halo density field using different methods: the grid-based method with Nearest Grid Point (NGP) and Cloud-In-Cell (CIC) assignment schemes and the Delaunay Tessellation (DT) method. We choose haloes above a limit between 10^{10} and $10^{11.5} h^{-1} M_{\odot}$ and reconstruct the smaller haloes using the conditional mass function of Equation (5). In this test, we assumed for $b(m)$ and $n(m)$ the forms calibrated on N-body simulations by Tinker et al. (2008) and Tinker et al. (2010). The output of the reconstruction is illustrated in Fig. 1, which shows the spatial distribution of original and reconstructed haloes in a thin slice of the Millennium simulation.

To test the accuracy of the method we perform the reconstruction on the MultiDark simulation, which gives us a better probe of the large-scale halo clustering. We measure the halo bias in the low-mass regime from the reconstructed halo catalogue. The halo bias has been estimated by first measuring the halo power spectrum $P(k)$ and then taking the square root of the ratio between the halo power spectrum and that of mass. In this, we assumed the non-

linear mass power spectrum given by CosmicEmu (Lawrence et al. 2010).

The recovered halo biases in mass bins below the resolution limit are shown in Fig. 2, which compares the results of using different estimates of the halo density field as well as different biasing models. In this figure, the measured halo bias is shown as a function of the wavenumber for the three mass bins: $10^{10} < m < 10^{10.5} h^{-1} M_{\odot}$, $10^{10.5} < m < 10^{11} h^{-1} M_{\odot}$, and $10^{11} < m < 10^{11.5} h^{-1} M_{\odot}$. We find that the DT method as implemented in the DTFE code (Cautun & van de Weygaert 2011) provides better results than the grid-based estimator with CIC and NGP assignment schemes. The large-scale bias, expected to asymptote to linear theory predictions, is in very good agreement with the predictions of Tinker et al. (2010) in the case of DT, whereas for the other methods the bias is clearly overestimated. This is particularly true in the case of NGP. The DT method better accounts for local variations in number density, reducing the shot noise in the reconstruction and giving a better sampling of the most extreme environments. In this exercise, we pushed the methods towards their limits by considering a very small grid size of $2.5 h^{-1} \text{ Mpc}$. However, if we increase the grid size to $5 - 10 h^{-1} \text{ Mpc}$, the recovered halo biases come to agreement and we find that the three methods converge to the same values.

The biasing scheme that enters in the conditional mass function has also some impact on the recovered halo clustering, in particular for small grid size density field reconstruction such as the one considered here. We show in the bottom panel of Fig. 2 the effect on the recovered halo bias when assuming a linear or power-law bias model as describe in Section 2.2. In both cases we use the halo density field reconstructed with the DT method. We find that

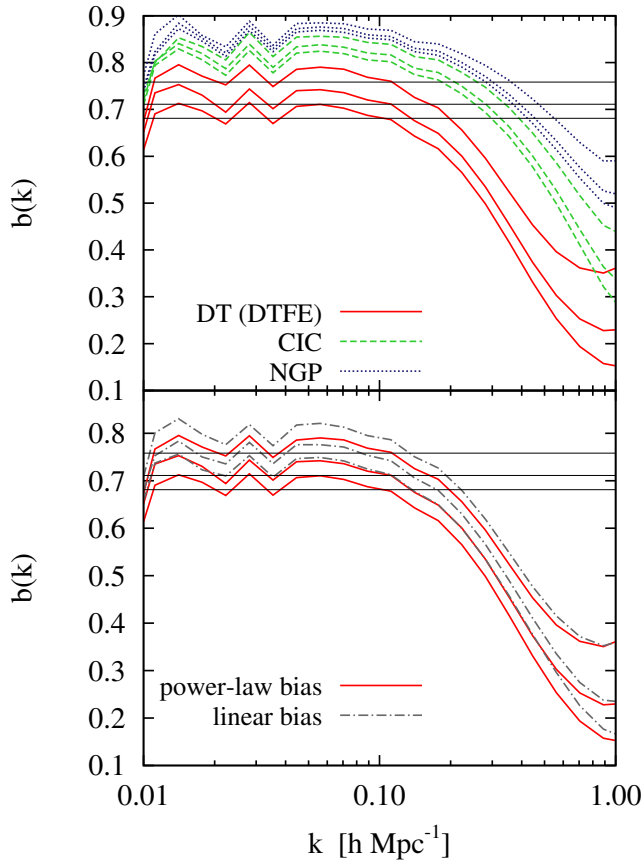


Figure 2. Top: scale-dependent bias of the reconstructed haloes obtained using different methods applied to the MultiDark data to estimate a continuous halo density field. For each of the methods: DT (dotted), CIC (dashed), and NGP (dot-dashed), the three curves correspond to samples of haloes with mass in the ranges $10^{10} < m < 10^{10.5} h^{-1} M_{\odot}$, $10^{10.5} < m < 10^{11} h^{-1} M_{\odot}$, and $10^{11} < m < 10^{11.5} h^{-1} M_{\odot}$, respectively from top to bottom, reconstructed using the haloes above the upper mass limit in each case. Bottom: scale-dependent bias of the reconstructed haloes obtained using linear and power-law bias prescriptions in the reconstruction. In these cases, the DT method has been used to estimate a continuous halo density field. In all panels, the horizontal lines are the linear halo bias predictions by Tinker et al. (2010) for the three mass bins considered.

the linear model tends to overestimate the large-scale linear bias for low-mass haloes compared to the power-law model, which instead allows us to recover the linear bias predictions of Tinker et al. (2010) at the few percent level.

It is noticeable in Fig. 2 that, as is inevitable, one cannot reconstruct the highest k regime of the halo power spectrum. This however does not really matter for the purpose of galaxy mock construction, since the overall galaxy power spectrum is dominated by the 1-halo term in this regime, as we will show in the next section.

4 APPLICATION TO GALAXY MOCK SAMPLE CONSTRUCTION

The reconstruction of the halo density field below the resolution limit in cosmological simulations is particularly useful in the context of building realistic galaxy mock surveys. As explained earlier, forthcoming large cosmological surveys of galaxies will need a large number of mock survey realisations, and we need these mocks

to include galaxies of very low luminosity/stellar mass. These dim galaxies sit in low-mass haloes, so that a method such as the present one is required to restore such missing haloes. In the following, we apply our halo reconstruction method and test its efficiency in this context.

An efficient way to build galaxy mock samples is to use the Halo Occupation Distribution (HOD) formalism (Seljak 2000; Peacock & Smith 2000; Cooray & Sheth 2002), which enable us to populate haloes with galaxy in a way to accurately reproduce the galaxy clustering. We use this technique on the Millennium simulation to build galaxy catalogues mimicking Sloan Digital Sky Survey DR7 (SDSS, Abazajian et al. 2009) volume-limited samples at $z \simeq 0.1$. We create two absolute magnitude-selected samples corresponding to $M_r - 5 \log(h) < -18$ and $M_r - 5 \log(h) < -19$ from the halo occupation measurements performed by Zehavi et al. (2011). We choose these cuts because they involve a significant fraction of the galaxies residing in the low-mass end of the halo mass function. In practice to create the galaxy catalogues, we populate haloes with central and satellite galaxies using the mean occupation numbers given by the HOD. While central galaxies are distributed at halo centres, satellite galaxies are randomly disposed around halo centres in such a way that their radial distribution follows a NFW (Navarro et al. 1996) density profile. The details of the procedure are given in Appendix B of de la Torre & Guzzo (2012). For each volume-limited sample we construct three catalogues: one based on the original complete halo catalogue to which we refer in the following as the fiducial sample; a second built from a reconstructed halo catalogue below $m_{\text{lim}} = 10^{11.5} h^{-1} M_{\odot}$ using $G = 2.5 h^{-1} \text{ Mpc}$; and a third one built from a reconstructed halo catalogue below $m_{\text{lim}} = 10^{11} h^{-1} M_{\odot}$ using $G = 1 h^{-1} \text{ Mpc}$. In these reconstructions, we estimated the halo density field using the DT method and assumed the power-law bias model in the conditional mass function.

We present in Fig. 3 the galaxy two-point correlation functions for the two absolute magnitude-selected mock samples built from the original complete halo catalogue. These are compared to those measured in the mock samples in which the haloes of mass below $m_{\text{lim}} = 10^{11.5} h^{-1} M_{\odot}$ and $m_{\text{lim}} = 10^{11} h^{-1} M_{\odot}$ have been reconstructed. In the case of the reconstruction with $m_{\text{lim}} = 10^{11.5} h^{-1} M_{\odot}$ and $G = 2.5 h^{-1} \text{ Mpc}$, we find that the correlation functions obtained are well recovered, although the correlation function is underestimated by up to 15% on intermediate scales for the $M_r - 5 \log(h) < -18$ sample. This underestimation is a direct consequence of the resolution scale chosen for the reconstruction. Indeed, the clustering drops on smaller scales than the reconstruction scale for the reconstructed low-mass haloes. This in turn causes the observed underestimation of the correlation function of less luminous galaxies. However, by reconstructing the halo density field on smaller scales, i.e. by using a lower mass limit for the reconstruction of $m_{\text{lim}} = 10^{11} h^{-1} M_{\odot}$, one can better reproduce the halo clustering on $1 h^{-1} \text{ Mpc}$ scales and eliminate the underestimation as shown in Fig. 3. We note that the small-scale galaxy clustering, i.e. below $0.7 - 1 h^{-1} \text{ Mpc}$, is always well recovered. This comes from the fact that the galaxy distribution inside haloes is independent of the halo clustering by construction (de la Torre & Guzzo 2012) and the clustering on those scales is dominated by this 1-halo term.

One could improve the method in the case of relatively coarse grid reconstructions, by working at the sub-grid level and using additional constraints from the mass non-linear power spectrum. Indeed, one could envisage distributing haloes in each sub-cell so as to reproduce the non-linear correlation function predicted from the-

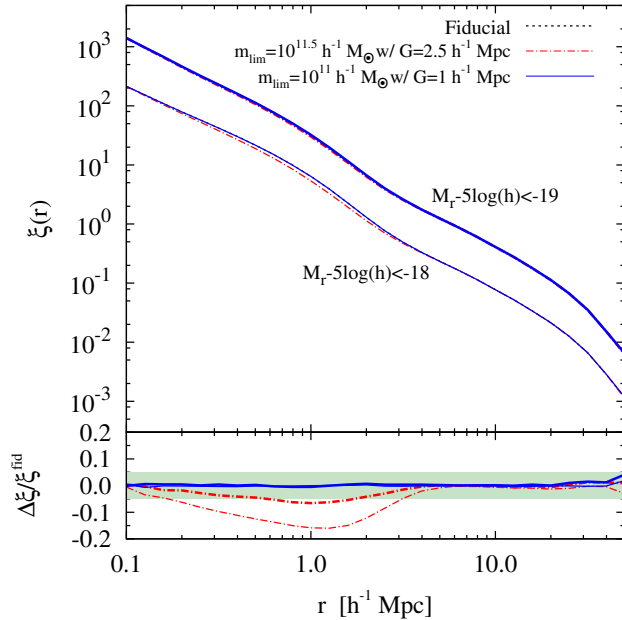


Figure 3. Top panel: comparison of the two-point correlation functions of mock galaxies with $M_r - 5 \log(h) < -18$ (thin lines) and $M_r - 5 \log(h) < -19$ (thick lines) obtained from the original halo catalogue (referred as fiducial in the caption) with those obtained after reconstruction of haloes below $10^{11.5} h^{-1} M_\odot$ (dot-dashed lines) and $10^{11} h^{-1} M_\odot$ (solid lines) in the Millennium simulation. The two reconstructions have been performed respectively on grid sizes of $G = 2.5 h^{-1} \text{ Mpc}$ and $G = 1 h^{-1} \text{ Mpc}$. The amplitude of the correlation functions for the $M_r - 5 \log(h) < -18$ samples have been divided by 5 to improve the clarity of the figure. Bottom panel: relative difference of two-point correlation functions of samples with reconstructed haloes with respect to that of fiducial samples. In the two panels the dotted and solid curves are not distinguishable as they almost overlap completely.

ory, instead of randomly distributing them. We plan to investigate this extension of the method elsewhere.

5 SUMMARY AND CONCLUSIONS

We have described in this letter a method for populating dark matter simulations with haloes of mass below the resolution limit. It is based on estimating a continuous halo density field and then sampling this stochastically in order to obtain low-mass haloes with the correct abundance and bias. This latter part requires the conditional halo mass function, which is extrapolated from the simulation itself or taken from theoretical predictions.

We found that the method works accurately and allows us to reproduce the halo distribution below the resolution limit with high fidelity, in particular for reconstructions on grids of size $G = 1 h^{-1} \text{ Mpc}$ or below. Moreover, the method is particularly efficient at producing galaxy mock samples from low-resolution simulation halo catalogues. We built galaxy mock samples using the HOD technique on the reconstructed halo density field, tuned to mimic SDSS observations. We showed that within a reasonable resolution limit range, one can recover the two-point correlation function at the percent level. Our results demonstrate that one can communicate efficiently the full information content of a large simulation by using only a catalogue of the more massive haloes. This method

should be very useful in the future in building realistic galaxy mocks for the massive forthcoming cosmological surveys such as EUCLID (Laureijs et al. 2011), where the volume and mass resolution requirements for the survey simulations are both very high.

It is also important to emphasise that the method presented here is relatively general. Another possible application is the construction of catalogues of dark matter particles. As with the galaxy mock sample construction one can make use of the halo model to distribute dark matter particles inside the haloes. Such catalogues could then be used to create cosmic shear catalogues from ray-tracing through the simulation, or to predict the galaxy-lensing signal, a quantity directly related to the galaxy-mass correlation function.

ACKNOWLEDGEMENTS

We thank Gabriella De Lucia for giving us useful comments on the manuscript. The MultiDark Database used in this paper and the web application providing online access to it were constructed as part of the activities of the German Astrophysical Virtual Observatory as result of a collaboration between the Leibniz-Institute for Astrophysics Potsdam (AIP) and the Spanish MultiDark Consolider Project CSD2009-00064. The Bolshoi and MultiDark simulations were run on the NASA's Pleiades supercomputer at the NASA Ames Research Center.

REFERENCES

- Abazajian K. N. et al., 2009, *ApJS*, 182, 543
- Bardeen J. M., Bond J. R., Kaiser N., Szalay A. S., 1986, *ApJ*, 304, 15
- Cautun M. C., van de Weygaert R., 2011, *ArXiv e-print 1105.0370*
- Cole S., Kaiser N., 1989, *MNRAS*, 237, 1127
- Cooray A., Sheth R., 2002, *Phys. Rep.*, 372, 1
- de la Torre S., Guzzo L., 2012, *MNRAS*, 427, 327
- Laureijs R. et al., 2011, *ArXiv e-print 1110.3193*
- Lawrence E., Heitmann K., White M., Higdon D., Wagner C., Habib S., Williams B., 2010, *ApJ*, 713, 1322
- Mann R. G., Peacock J. A., Heavens A. F., 1998, *MNRAS*, 293, 209
- Narayanan V. K., Berlind A. A., Weinberg D. H., 2000, *ApJ*, 528, 1
- Navarro J. F., Frenk C. S., White S. D. M., 1996, *ApJ*, 462, 563
- Peacock J. A., Smith R. E., 2000, *MNRAS*, 318, 1144
- Prada F., Klypin A. A., Cuesta A. J., Betancort-Rijo J. E., Primack J., 2012, *MNRAS*, 423, 3018
- Schlegel D. et al., 2011, *ArXiv e-print 1106.1706*
- Seljak U., 2000, *MNRAS*, 318, 203
- Springel V. et al., 2005, *Nature*, 435, 629
- Tinker J., Kravtsov A. V., Klypin A., Abazajian K., Warren M., Yepes G., Gottlöber S., Holz D. E., 2008, *ApJ*, 688, 709
- Tinker J. L., Robertson B. E., Kravtsov A. V., Klypin A., Warren M. S., Yepes G., Gottlöber S., 2010, *ApJ*, 724, 878
- Wild V. et al., 2005, *MNRAS*, 356, 247
- Zehavi I. et al., 2011, *ApJ*, 736, 59

This paper has been typeset from a \LaTeX file prepared by the author.